

# CMPE 561

## Natural Language Processing

### Application Project 2

Language Identification System

*Ezgi Tekdemir*  
*Beyza Gul Gurbuz*

submitted to  
Tunga GUNGOR

December 19, 2017

## Introduction

For this project, we implemented a language identification system in Python using Discriminating between Similar Languages (DSL) Shared Task 2015 corpus. The analysis is based on character unigrams, meaning that the identification of a language is based on letter frequencies. The hypothesis is that each language either has unique characters or use the characters with unique probabilities.

In our corpus, there are 13 languages, each with 2000 sentences, so our corpus has a balanced distribution with 26000 sentences in total. The languages in question are Bulgarian, Bosnian, Czech, Argentine Spanish, Peninsular Spanish, Croatian, Indonesian, Macedonian, Malay, Brazilian Portuguese, European Portuguese, Slovak and Serbian.

## Methods

We mainly tried two models to solve this problem.

### Generative Modeling (Naive Bayes)

In this approach, we first train our model by calculating the letter probabilities for each language using Laplace Smoothing.

$$P(c_i|l) = \frac{\text{number of times } c_i \text{ appears in } l + 1}{\text{total number of characters in } l + \text{number of unique characters in } l} \quad (1)$$

Then, for a sentence in the test set, we sum up the letter probabilities for each letter in the sentence and then multiply it with probability of that language in the corpus. Since there are same amount of sentences for each language, the probability for the language  $l$  is the same for all languages,  $2000/26000 = 1/13$ .

$$P(l|s) = \left( \sum_{i=1}^n P(c_i|l) \right) \cdot \frac{1}{13} \quad (2)$$

For a sentence, we do this for each language and say that the language with the highest probability is the predicted language for that sentence.

## Discriminative Modeling (SVM)

For this approach, we used Cornell’s SVM multiclass library. We used distinct letters that appear in the corpus as features and the existence of a character in that sentence as the corresponding value of that feature, i.e. the value of a feature is either 0 or 1.

While training, the system learns the important letters for each language. It uses a letter pool consisting of letters of all languages in the training set. Then, using this model, the system predicts the most likely language for that sentence.

As bonus, we added three other features to our SVM: count of letters, bigram letter frequencies and count of capital letters per sentence. The accuracy significantly increased with the first two features; however, the capital letter count decreased the accuracy of our predictions, so we decided to exclude it from the analysis at the end. We called the SVM with existence and count of unigram letter frequencies and count of bigram letter frequencies as *Super SVM*.

## Evaluation

We evaluated the approaches used above with three types of metrics.

### Accuracy

$$Accuracy = \frac{\text{number of correct predictions}}{\text{number of sentences}} \quad (3)$$

We calculated the accuracy for each language and for the entire test set.

### Averaged Measures (Precision, Recall, F-measure)

In order to calculate these measures, the following values must be calculated first:

**TP** True Positive: Number of sentences whose language is L and predicted as L

**FN** False Negative: Number of sentences whose language is L but predicted as another language

**FP** False Positive: Number of sentences whose language is another language but predicted as L

**Precision** This measure answers the question of "Of all sentences that are labeled with language l, how many actually belong to language l?". It can be calculated as  $\frac{TP}{TP+FP}$ .

**Recall** This measure answers the question of "Of all sentences that truly belong to language l, how many did we label as language l?". It can be calculated as  $\frac{TP}{TP+FN}$ .

**F-measure** This is the weighted average of precision and recall. This is a better measure than accuracy if the cost of false positives and false negatives are very different. It can be calculated as  $\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$ .

**Microaveraged Measures** Gives equal weight to sentence language pairs, so F-measure uses averaged precision and recall values.

**Macroaveraged Measures** Gives equal weight to languages, so takes the average of precision and recall values regardless of sentence count.

At the end, we compare our approaches using these performance measures and find out which approach and features work best to identify the languages of sentences.

## Program Interface

In order to run the program, Python3.5 must be installed. Put identifier.py in the same directory with the [corpus file](#) and rename the corpus file as "corpus.txt".

For the SVM part, download the source code from [Cornell's SVM-Multiclass website](#). Then, go to svm-multiclass directory from the terminal and type:

```
> make
```

Wait for the code to compile. The executables svm-multiclass-learn and svm-multiclass-classify will be produced, which will be used in our program.

The program can be run in three modes: naive Bayes, unigram SVM and super SVM. Note that the program will not run without the mode argument.

In order to run the program in naive bayes mode, go to the directory of the python file and type:

```
> python identifier.py naive\_bayes
```

In order to run the program in SVM modes, path to the svm-multiclass directory must also be given as an argument.

For unigram SVM mode, type:

```
> python identifier.py unigram\_svm path/to/svm-multiclass
```

For super SVM mode, type:

```
> python identifier.py super\_svm path/to/svm-multiclass
```

## Program Execution

It is required to give a mode to the program since there are three approaches that we use (naive bayes, unigram SVM and SVM with more features) to train our program.

For the naive bayes mode, The SVM modes also require the path to the svm-multiclass directory to run since the svm-multiclass-learn and svm-multiclass-classify executables used for training and testing the sentences are run inside our program using Python's *os* library. The commands executed are as follows:

```
> path/to/svm-multiclass/svm_multiclass_learn -c 1.0
    svm_train.txt svm_model.txt
> path/to/svm-multiclass/svm_multiclass_classify
    svm_test.txt svm_model.txt svm_output.txt
```

The first command is for training the system using the svm\_train file we created in our program. This command produces svm\_model file (in the directory of the Python file), which will be used in the next command to classify the test input. We chose the argument *c* (trade-off between training error and margin) to be 1.0 since we found out that the accuracy increases when this argument is higher. We do not allow the end-user to change this value.

The second command is for testing the system using the svm\_test file we created in our program. This command uses the svm\_model file and produces a file called svm\_output. In this file, there is one line per a line in the svm\_test file in the same order. The first value of each file is the predicted language number, and the following 13 numbers are the discriminant values for each of the language classes. In our program, we compare the first numbers of each lines of svm\_test and svm\_output to calculate the performance measures.

For all three modes, the produced output is of the same format. Accuracy, micro-averaged measures (Precision, Recall, F-measure) and macro-averaged measures (Precision, Recall, F-measure) are outputted to the terminal. In SVM modes, the multiclass library also produces an output before our output.

An example execution with unigram SVM mode is shown below:

```
Overall Accuracy: 43.07692307692308

Accuracy per Language:
pt-BR: 9.836065573770492
cz: 99.5
bs: 0.0
es-AR: 17.117117117117118
pt-PT: 35.5
sk: 6.310679611650485
sr: 0.0
hr: 55.55555555555556
bg: 99.05660377358491
mk: 100.0
my: 0.0
id: 88.71794871794872
es-ES: 46.63677130044843

Micro averaged Precision: 0.4307692307692308
Micro averaged Recall: 0.4307692307692308
Micro averaged F-measure: 0.43076923076923074

Macro averaged Precision: 0.48180332869146836
Macro averaged Recall: 0.4294082628077506
Macro averaged F-measure: 0.382090673270263
```

Figure 1: Example execution of naive bayes mode

## Input and Output

As the input, DSL Shared Task 2015 corpus is used. In this file, there are 26000 sentences with their language labels each separated with a new line character. A sentence and its language are separated with multiple space characters. The labels are translated as follows:

**bg** Bulgarian  
**bs** Bosnian  
**cz** Czech  
**es-AR** Argentine Spanish  
**es-ES** Peninsular Spanish  
**hr** Croatian  
**id** Indonesian  
**mk** Macedonian  
**my** Malay  
**pt-BR** Brazilian Portuguese  
**pt-PT** European Portuguese  
**sk** Slovak  
**sr** Serbian

Балканите са доста сеизмично активен регион и няма нищо чудно, че след „Фукушима“ българите са се замислили за допълнителни мерки за сигурност“, коментира пред московския електронен вестник „Газета“ руски банков анализатор. **bg**

Figure 2: An example Bulgarian sentence

The algorithms are applied on this corpus, by randomly separating the sentences as 90% training and 10% testing. Naive Bayes algorithm directly gives an output on the terminal. The SVM algorithms, however, require `svm.train` and `svm.test` files to run. These two files have the same structure. On each line, there is the language-id, features and values corresponding to



each sentence. In our case, there are 13 language-ids. In unigram SVM mode, there are 264 possible values since the unique character size summed up for all languages is 264.

language-id f1:v1 f2:v2 f3:v3...

For example, if the language-id for English is 2, character type is a feature and the count of those characters are the corresponding values, and the English alphabet is mapped to values from 1 to 26, the word "bee" would be represented with the following line:

2 2:1 5:2

```

1 12 1:1 2:1 3:1 4:1 5:1 6:1 7:1 8:1 9:1 10:1 11:1 12:1 13:1 14:1 15:1 16:1 17:1 18:1 19:1 20:1 21:1 22:1 23:1 24:1 25:1 26:1 27:1 28:1 29:1 30:1 31:1 32:1 33:1 34:1 35:1 36:1
37:1 38:1 39:1 40:1 41:1 42:1 43:1 44:1 45:1 46:1 47:1 48:1 49:1 50:1 51:1 52:1 53:1 54:1 55:1 56:1 57:1 58:1 59:1 60:1 61:1 62:1 63:1 64:1 65:1 66:1 67:1 68:1 69:1 70:1 71:1 72:1 73:1 74:1 75:1 76:1 77:1 78:1 79:1 80:1
294:2 295:1 296:3 297:9 298:3 299:1 300:2 301:4 302:4 303:1 533:1 534:1 535:4 536:1 537:2 538:1 539:2 540:1 541:1 542:1 543:1 544:1 545:1 546:1 547:2 548:1 549:1 550:1 551:1
552:1 553:1 554:1 555:1 556:2 557:3 558:1 559:1 560:1 561:1 562:2 563:1 564:2 565:2 566:4 567:1 568:1 569:4 570:1 571:1 572:2 573:3 574:2 575:4 576:1 577:1 578:1 579:1 580:1
581:1 582:1 583:1 584:1 585:1 586:1 587:1 588:1
2 12 1:1 2:1 3:1 4:1 5:1 6:1 7:1 8:1 9:1 10:1 11:1 12:1 13:1 14:1 15:1 16:1 17:1 18:1 19:1 20:1 21:1 22:1 23:1 24:1 25:1 26:1 27:1 28:1 29:1 30:1 31:1 32:1 33:1 34:1 35:1 36:1 37:1 38:1 39:1 40:1 41:1 42:1 43:1 44:1 45:1 46:1
47:1 48:1 49:1 50:1 51:1 52:1 53:1 54:1 55:1 56:1 57:1 58:1 59:1 60:1 61:1 62:1 63:1 64:1 65:1 66:1 67:1 68:1 69:1 70:1 71:1 72:1 73:1 74:1 75:1 76:1 77:1 78:1 79:1 80:1
302:10 303:1 304:1 305:6 306:1 307:3 308:3 309:1 310:3 311:3 312:3 313:1 314:2 315:1 316:1 533:2 534:2 535:2 536:2 537:2 538:1 539:4 570:1 574:2 575:1 580:2 582:1
583:2 586:2 589:1 590:1 591:1 592:1 593:1 594:1 595:1 596:1 597:3 598:2 599:3 600:1 601:1 602:1 603:1 604:1 605:1 606:1 607:1 608:1 609:1 610:2 611:4 612:1 613:1 614:1 615:2
616:2 617:2 618:2 619:1 620:1 621:1 622:1 623:1 624:4 625:1 626:3 627:1 628:1 629:2 630:1 631:2 632:1 633:2 634:2 635:1 636:4 637:2 638:1 639:2 640:1 641:1 642:1 643:1 644:2
645:1 646:1 647:1 648:2 649:1 650:1 651:1 652:1 653:1 654:1 655:1 656:2 657:1 658:1
3 13 2:1 3:1 4:1 5:1 6:1 7:1 8:1 9:1 10:1 11:1 12:1 13:1 14:1 15:1 16:1 17:1 18:1 19:1 20:1 21:1 22:1 23:1 24:1 25:1 26:1 27:1 28:1 29:1 30:1 31:1 32:1 33:1 34:1 35:1 36:1 37:1 38:1 39:1 40:1 41:1 42:1 43:1 44:1 45:1 46:1
47:1 48:1 49:1 50:1 51:1 52:1 53:1 54:1 55:1 56:1 57:1 58:1 59:1 60:1 61:1 62:1 63:1 64:1 65:1 66:1 67:1 68:1 69:1 70:1 71:1 72:1 73:1 74:1 75:1 76:1 77:1 78:1 79:1 80:1
292:14 293:7 296:6 297:3:1 298:6 301:10 302:5 303:1 311:2 315:4 317:4 318:4 534:2 555:2 556:1 560:1 566:4 575:2 586:2 600:1 613:1 618:1 627:2 632:1 642:13 643:4 645:2 647:1
659:1 660:2 661:1 662:4 663:2 664:3 665:1 666:1 667:4 668:2 669:1 670:2 671:1 672:1 673:3 674:1 675:1 676:1 677:1 678:1 679:1 680:2 681:4 682:4 683:1 684:2 685:4 686:3 687:2
688:2 689:1 690:1 691:1 692:2 693:1 694:1 695:3 696:2 697:2 698:1
4 9 2:1 3:1 4:1 5:1 6:1 7:1 8:1 9:1 10:1 11:1 12:1 13:1 14:1 15:1 16:1 17:1 18:1 19:1 20:1 21:1 22:1 23:1 24:1 25:1 26:1 27:1 28:1 29:1 30:1 31:1 32:1 33:1 34:1 35:1 36:1 37:1 38:1 39:1 40:1 41:1 42:1 43:1 44:1 45:1 46:1 47:1 48:1 49:1 50:1 51:1 52:1 53:1 54:1 55:1 56:1 57:1 58:1 59:1
57:13 572:1 576:1 578:1 579:2 580:3 585:1 586:4 593:1 597:2 598:2 600:2 604:2 608:2 624:2 635:3 636:1 642:4 643:3 647:1 651:1 655:1 671:3 673:3 675:1 682:2 684:1 686:1 697:2
699:1 700:1 701:3 702:1 703:1 704:1 705:2 706:4 707:1 708:2 709:2 710:1 711:1 712:4 713:1 714:2 715:2 716:3 717:1 718:1 719:1 720:4 721:2 722:2 723:3 724:1 725:3 726:1 727:2
728:2 729:1 730:1 731:1 732:1 733:1 734:1 735:1 736:2 737:1 738:4 739:1 740:1
5 4 1:1 2:1 3:1 4:1 5:1 6:1 7:1 8:1 9:1 10:1 11:1 12:1 13:1 14:1 15:1 16:1 17:1 18:1 19:1 20:1 21:1 22:1 23:1 24:1 25:1 26:1 27:1 28:1 29:1 30:1 31:1 32:1 33:1 34:1 35:1 36:1 37:1 38:1 39:1 40:1 41:1 42:1 43:1 44:1 45:1 46:1 47:1 48:1 49:1 50:1 51:1 52:1 53:1 54:1 55:1 56:1 57:1 58:1 59:1
60:1 61:1 62:1 63:1 64:1 65:1 66:1 67:1 68:1 69:1 70:1 71:1 72:1 73:1 74:1 75:1 76:1 77:1 78:1 79:1 80:1
312:1 314:3 315:4 317:3 320:5 321:2 322:2 323:2 324:2 325:1 326:1 327:2 334:1 547:3 555:2 557:3 562:2 565:1 566:1 568:1 569:3 570:1 574:1 580:4 583:2 599:3 600:1 608:1 624:1
630:1 631:3 639:1 649:2 654:1 668:1 663:2 668:1 679:1 688:2 694:4 704:1 715:1 717:3 720:1 725:2 741:1 742:2 743:3 744:2 745:1 746:1 747:1 748:2 749:1 750:1 751:1 752:1 753:1
754:1 755:1 756:2 757:1 758:1 759:1 760:1 761:1 762:1 763:2 764:1 765:2 766:1 767:1 768:3 769:1 770:1 771:1 772:1 773:2 774:1 775:1 776:1 777:1 778:1 779:1 780:1 781:1 782:2
783:1 784:1 785:1 786:2 787:1
6 5 2:1 3:1 4:1 5:1 6:1 7:1 8:1 9:1 10:1 11:1 12:1 13:1 14:1 15:1 16:1 17:1 18:1 19:1 20:1 21:1 22:1 23:1 24:1 25:1 26:1 27:1 28:1 29:1 30:1 31:1 32:1 33:1 34:1 35:1 36:1 37:1 38:1 39:1 40:1 41:1 42:1 43:1 44:1 45:1 46:1 47:1 48:1 49:1 50:1 51:1 52:1 53:1 54:1 55:1 56:1 57:1 58:1 59:1
270:15 271:8 272:10 274:1 279:1 280:1 283:5 284:11 285:1 286:17 289:1 293:7 296:2 297:19 298:8 299:1 300:10 301:7 302:9 303:2 305:2 311:3 315:3 316:3 320:1 328:2 329:1 330:1
331:1 332:1 333:1 336:2 347:3 555:2 560:3 562:2 566:2 574:1 580:1 583:4 585:2 599:3 604:1 611:1 613:1 618:1 624:3 626:4 629:2 641:1 642:3 645:2 647:1 649:1 652:5 660:3 664:2
666:2 669:1 675:2 682:1 683:2 688:4 702:1 706:1 714:7 717:2 733:2 761:2 788:1 789:1 790:3 791:1 792:1 793:3 794:1 795:1 796:1 797:1 798:2 799:1 800:1 801:1 802:1 803:1 804:1
805:2 806:2 807:5 808:1 809:1 810:1 811:1 812:2 813:1 814:1 815:1 816:1 817:1 818:1 819:1 820:1 821:1 822:1 823:1 824:1
7 3 2:1 3:1 4:1 5:1 6:1 7:1 8:1 9:1 10:1 11:1 12:1 13:1 14:1 15:1 16:1 17:1 18:1 19:1 20:1 21:1 22:1 23:1 24:1 25:1 26:1 27:1 28:1 29:1 30:1 31:1 32:1 33:1 34:1 35:1 36:1 37:1 38:1 39:1 40:1 41:1 42:1 43:1 44:1 45:1 46:1 47:1 48:1 49:1 50:1 51:1 52:1 53:1 54:1 55:1 56:1 57:1 58:1 59:1
273:1 283:3 284:9 286:6 289:1 290:1 292:9 293:4 296:1 297:30 298:8 299:1 300:2 301:5 302:12 303:2 305:1 311:4 315:9 316:5 317:4 318:2 325:1 534:2 536:3 555:6 560:3 565:1
566:1 571:1 575:3 579:1 586:3 590:2 608:3 611:2 626:5 627:3 630:1 632:3 635:1 642:5 643:1 649:3 650:1 653:1 660:2 662:3 664:3 667:2 670:2 671:2 673:2 675:1 688:2 691:2 694:3
695:3 723:2 768:1 709:1 807:3 812:2 823:2 824:1 825:1 826:3 827:1 828:1 829:1 830:1 831:1 832:1 833:2 834:1 835:1 836:1 837:1 838:1 839:7 840:1 841:1 842:1 843:1 844:1 845:1
4 2:1 3:1 4:1 5:1 6:1 7:1 8:1 9:1 10:1 11:1 12:1 13:1 14:1 15:1 16:1 17:1 18:1 19:1 20:1 21:1 22:1 23:1 24:1 25:1 26:1 27:1 28:1 29:1 30:1 31:1 32:1 33:1 34:1 35:1 36:1 37:1 38:1 39:1 40:1 41:1 42:1 43:1 44:1 45:1 46:1 47:1 48:1 49:1 50:1 51:1 52:1 53:1 54:1 55:1 56:1 57:1 58:1 59:1
852:1 853:1 854:1 855:1 856:1 857:1 858:1 859:2 860:2 861:1 862:1 863:1 864:1 865:1 866:1 867:2 868:1 869:1 870:1 871:1 872:1 873:1
8 8 1:1 2:1 3:1 4:1 5:1 6:1 7:1 8:1 9:1 10:1 11:1 12:1 13:1 14:1 15:1 16:1 17:1 18:1 19:1 20:1 21:1 22:1 23:1 24:1 25:1 26:1 27:1 28:1 29:1 30:1 31:1 32:1 33:1 34:1 35:1 36:1 37:1 38:1 39:1 40:1 41:1 42:1 43:1 44:1 45:1 46:1 47:1 48:1 49:1 50:1
59:1 69:1 70:1 71:1 72:1 267:1 268:8 269:19 270:12 271:25 272:8 273:22 280:1 283:9 284:10 285:1 286:10 289:2 291:1 292:9 293:8 294:4 295:1 296:6 297:22 298:5 299:1 300:1
301:6 302:6 303:2 305:3 306:1 307:3 308:9 309:1 311:5 313:1 314:5 316:3 325:1 335:1 336:1 337:1 338:3 537:2 536:2 560:2 562:2 565:1 566:3 568:2 569:4 570:1 571:1 575:2 579:1
580:3 582:2 583:3 586:1 590:4 591:1 593:1 596:1 597:4 598:2 600:2 604:2 622:1 623:3 624:4 625:2 626:2 627:3 631:3 633:2 636:4 641:1 651:1 660:1 670:2 674:1 675:1 683:1 688:1
691:1 694:2 697:2 704:4 705:1 709:1 716:1 717:2 720:3 722:2 728:1 733:1 734:1 736:1 743:1 763:1 767:1 768:1 783:1 785:2 789:2 797:3 823:1 837:2 840:1 871:1 874:1 875:2 876:1
877:1 878:2 879:1 880:1 881:4 882:2 883:1 884:1 885:1 886:1 887:2 888:1 889:1 890:1 891:1 892:1 893:1 894:1 895:1 896:1 897:1 898:1 899:1 900:1 901:1 902:3 903:1 904:1 905:1

```

Figure 3: Example train file (same format with the test file)

In the svm\_model file, the learned rule from the training using the regularization parameter c (in our case, 1.0) is stored.

```

SVM-multiclass Version V2.20
13 # number of classes
11823 # number of base features
0 # loss function
0 # kernel type
3 # kernel parameter -d
1 # kernel parameter -g
1 # kernel parameter -s
1 # kernel parameter -r
empty# kernel parameter -u
143300 # highest feature index
61 # number of training documents
2 # number of support vectors plus 1
0 # threshold b, each following line is a SV (starting with alpha)
1 qid:0 1:-0.001253346 2:-0.017148096 3:-0.017694868 4:-0.017687812 5:-0.017672254 6:-0.0033349127 7:-0.017398862 8:-0.0022987763 9:-0.0059565897 10:-0.00047999344
11:-0.00066202815 12:5.673429e-05 13:-0.0004495516 14:-0.00012196952 15:-0.0019465407 16:-0.0002520242 17:-0.017282664 18:-0.017500903 19:-0.010586752 20:-0.017638077
21:0.0003595096 22:-0.003828626 23:-0.0061863577 24:-0.003072453 25:-0.0018349669 26:-0.0087951487 27:-0.017656485 28:-0.0043193693 29:0.0026036517 30:-0.011292303
31:-0.017698519 32:-0.017654184 33:-0.0035185787 34:-0.016160684 35:-0.017577443 36:-0.01770461 37:-0.017731722 38:-0.0046796849 39:-0.015560844 40:-0.0038408625
41:-0.00012141965 42:-0.012592169 43:-0.0027638006 44:-0.00043806501 45:-0.012271672 46:-0.016830032 47:0.0038261695 48:-0.0027448403 49:-0.0086659947 50:-0.0037689924
51:0.00069276184 52:-0.0038376446 53:-0.0003666118 54:-0.003082375 55:-0.00030931653 56:0.0033751687 57:-0.00024592809 58:-0.0001307344 59:-0.0028435937 60:-1.743597e-05
61:-0.002232854 62:0.0018369895 63:-0.0052699526 64:-0.006168785 65:-0.0023561991 66:-0.00064873043 67:-0.0021109416 68:-0.0045033363 69:-0.001189078 70:-0.000123669
71:-7.4622825e-05 72:-0.0007436529 73:-0.0096347658 74:-0.010804166 75:-0.0056373398 76:-0.0042182626 77:-0.0017882166 78:-0.0033546942 79:-0.003746932 80:-0.00065294735
81:-0.0005632298 82:-2.9266927e-05 83:-0.0005595746 84:-0.000657368 85:-0.00026376337 86:-0.0013369725 87:-3.4167253e-06 88:-0.0019316794 89:-0.000211169 91:-0.00013077108
92:-0.0030794619 93:-0.0030864532 94:-0.0001874418 95:-0.00046577878 96:-0.0001970401 97:-0.00088384602 98:-0.00088384602 99:-0.00088384602 101:-0.00088384602
102:-0.00088384602 103:-0.0002355964 104:-0.00088384602 105:-0.00088384602 106:-0.00088384602 107:-0.00082738304 108:-0.00088384602 109:-0.00086480711 110:-0.00088384602
111:-0.00086498071 112:-0.0008350802 113:-0.00037419074 114:-0.00017486753 115:-0.00016054831 116:-0.00085393345 118:-0.00011292232 119:-0.00079755409 121:-0.00066653674
122:-0.00018634058 123:-0.00061785791 124:-0.0007639166 126:-0.00083362725 127:-0.0036780851 128:-0.00070866293 129:0.001655492 130:0.013492111 131:0.007995952
132:0.013472541 133:-0.0001438231 134:-0.015585664 135:0.0055064149 136:-0.00033187261 137:-6.524182e-06 138:-0.0001725663 139:-0.00038994756 140:0.0002893921
141:-0.0011804207 142:-0.0008073432 143:-0.0021491961 144:-0.00028121116 145:-0.00086289071 146:-0.00021596934 147:-0.000125437 148:-0.00031849722 149:-0.00089817892
150:-0.0001861134 151:-1.0873638e-05 152:-9.9769486e-05 153:-0.0001371122 155:-0.00014066089 157:-0.00018026333 158:-0.0019162615 159:-7.6898825e-05 160:-0.00018526199
161:-0.00022691539 162:-0.0002186297 163:-0.0002170572 164:-0.00011649238 165:-0.0002342167 166:0.0004459741 167:-0.00082804257 168:-4.7243921e-05 169:-0.0001776827
170:-9.4998086e-05 171:-0.00039407899 172:-0.00028127237 173:-6.582285e-05 174:-0.0016237866 175:-0.00070254371 176:-7.1491915e-05 177:-5.0669674e-05 178:-6.0837865e-05
179:-5.0669674e-05 180:-2.1747276e-06 181:-7.1656999e-05 182:5.4855609e-05 183:-0.00033801439 184:-0.00014453036 185:0.00012820512 186:-5.2138956e-05 187:-2.6096728e-05
188:-5.2138956e-05 189:-2.1282243e-06 190:-0.0001486295 191:-0.00045818402 192:-6.524182e-06 193:-2.1282243e-06 194:-4.7387807e-05 195:-6.524182e-06 196:2.4544106e-05
197:-1.0071638e-05 199:-0.0003234187 201:-5.123135e-05 202:-6.689183e-06 204:-0.0005302e-05 205:0.00013673868 207:-3.4686634e-05 208:-1.2993862e-05 209:-3.4167253e-06
210:-0.00018789296 211:-0.0011404781 212:-1.1938422e-05 214:-2.3791001e-06 215:-0.00025073742 222:-3.5208894e-15 223:0.00012396468 224:-1.1938422e-05 225:-2.1282243e-06
229:-4.3494551e-06 233:4.2735042e-05 234:4.2735042e-05 235:-2.1747276e-06 237:4.2735042e-05 238:-2.1747276e-06 241:-2.0280857e-05 243:-4.2735042e-05 252:-4.2735042e-05

```

Figure 4: Example model file

In addition to terminal output, SVM algorithms give an output file called `svm_output.txt`. In the `svm_output` file, there are 14 numbers per line. The first number denotes the language id and the other 13 are the discriminant values for each of the language classes (values of  $x \cdot w_i$ ).

```

1 1 1.058748 0.513624 -0.343680 -0.289460 0.535259 -0.271806 -0.927521 0.439519 -0.125423 0.867802 0.046700 0.164319 -1.668080
2 0 0.970500 0.144677 0.387373 -0.052790 0.707876 -0.211918 0.509613 0.274895 0.255897 0.445447 0.022077 0.283738 -3.737385
3 1 3.086361 0.339907 -0.460615 -0.283309 1.317885 -0.312253 -0.291336 0.415400 -0.699486 0.557482 -0.194305 0.213764 -4.161965
4 1 5.427884 0.615121 -0.536092 -0.732228 2.079143 -0.470376 -3.535466 1.248478 0.081740 1.431100 -0.336035 0.818377 -6.091645
5 1 2.446384 0.591722 0.273131 0.239190 0.957831 -0.411390 -0.827424 0.416147 -0.129849 0.180308 0.016111 0.173344 -3.925423
6 2 -0.046460 -1.372880 -0.908590 -1.190624 -1.305525 15.011465 -1.070947 -0.754543 -1.129117 -2.510448 0.175854 -0.866821 -2.031364
7 1 6.982110 0.051128 -0.641796 -0.284969 1.886108 -0.624168 -0.354003 0.441763 -0.826537 0.262425 -0.205284 -0.311152 -6.377444
8 13 0.539952 -0.632712 0.974619 -1.891523 -1.459629 -0.307100 -2.634197 0.678283 -0.461676 -1.744802 -0.219517 -0.282573 7.440877
9 1 3.211593 -0.174827 -0.388985 -0.224002 0.703785 -0.193468 -1.709679 0.633544 -0.045981 0.154416 -0.194152 0.089348 -1.861591
10 1 1.187408 0.875977 -0.132455 0.354725 0.603316 -0.202952 -0.270704 0.372595 0.018535 0.788137 -0.044437 0.131416 -3.681556
11 7 0.536807 -0.447892 -0.544095 -0.305357 -0.411441 -0.199837 5.248810 -0.358790 -0.604080 -0.781537 -0.054165 -0.508435 -1.569988
12 13 -0.339953 0.037010 0.650171 -0.265436 -0.619453 -0.256342 -1.287139 0.203086 -0.315439 -0.432882 0.036374 -0.180827 2.769929
13 1 2.972568 0.833958 0.742257 -0.478758 1.551686 -0.422608 -0.070436 0.797806 0.070735 0.675821 -0.087673 0.369711 -6.955065
14 8 0.126860 0.338263 -0.287826 0.239593 -0.056388 -0.152388 -1.493063 1.164367 0.715411 -0.069309 -0.081206 0.872598 -1.063192
15 0 0.997820 0.040580 -0.001416 -0.029371 0.470262 -0.079866 0.446227 0.182377 -0.183965 0.250812 -0.045743 -0.094703 -1.953012
16 10 1.140321 1.149601 -0.342041 -0.082384 0.700001 -0.192935 0.529145 0.546878 -0.039823 1.361333 -0.107032 0.279500 -4.942563
17 3 -0.899070 0.340311 0.961493 0.893189 -0.240984 -0.080010 -0.612669 0.849670 0.752039 0.182624 -0.089430 0.612697 -2.669859
18 1 1.922431 0.544695 -0.330585 0.515755 0.454413 -0.335006 -1.258339 0.404009 -0.142353 0.031832 -0.065575 0.080216 -1.821492
19 4 -0.709777 -0.019122 -1.323385 3.266811 -0.782695 -0.187421 2.351431 0.123602 0.215948 -0.287083 -0.059961 0.194607 -2.782957
20 13 -0.490941 0.144880 0.875674 -0.966258 -0.580012 -0.131711 -1.382337 0.346568 -0.194365 -0.378938 -0.141477 0.022449 2.876469
21 9 -0.737940 0.244425 0.238789 0.206797 -0.191082 -0.121112 -0.706594 0.834344 1.106539 -0.151736 -0.091001 0.812683 -1.443934
22 2 -0.211260 1.311312 -0.006612 0.364959 0.235416 -0.140435 -2.759553 0.356290 0.297022 0.510357 -0.021672 0.413010 -0.348834
23 7 -0.073769 0.047666 -0.270783 -0.394111 -0.219322 -0.140280 3.148876 0.064358 -0.204150 -0.000205 -0.038825 -0.162104 -1.757352
24 9 -0.507726 0.391505 -0.525688 0.394936 -0.403565 -0.152637 -1.200666 0.823988 1.040513 -0.297300 -0.065776 0.772457 -0.701800
25 8 -0.960210 0.586274 -0.694375 0.414117 -0.389304 -0.196331 -0.293192 1.716922 1.128033 0.086552 -0.189068 1.202110 -2.411527
26 4 -0.051298 1.278124 -0.231562 1.538724 0.336049 -0.162492 -1.918144 0.379856 0.350053 0.733039 -0.027524 0.413221 -2.638046
27 6 -1.654100 -1.127785 -0.718330 -0.906036 -1.078393 7.343685 -0.865886 -0.654597 -0.905284 -2.043804 5.033987 -0.682030 -1.661348
28 1 1.157840 0.308436 0.149041 -0.211876 0.858694 -0.342383 -0.128177 -0.000050 -0.056727 0.675681 0.080539 0.172783 -2.663000
29 6 -2.048747 -1.474155 -0.959299 -1.208125 -1.435200 17.026497 -1.182861 -0.828464 -1.166161 -2.625998 -1.057946 -0.855834 -2.093707
30 1 2.512766 0.527390 -0.741557 0.002135 0.797532 -0.235353 -1.268173 0.320331 -0.641914 0.624015 -0.131727 -0.321771 -1.443676
31 6 -1.593631 -1.076447 -0.736691 -0.934637 -1.005330 9.698203 -0.848751 -0.577070 -0.800356 -1.960239 2.102019 -0.660373 -1.526696
32 8 -0.597048 1.051012 0.081846 0.794518 -0.207024 -0.275291 -2.442294 1.222281 0.952883 0.172871 0.034541 0.942640 -1.730854

```

Figure 5: Example output file

We use `svm_output.txt` and `svm_test.txt` as inputs to our program again to compare the predicted and correct languages and evaluate our results. We only use the language-ids to compare these files.

Finally, as the terminal output, we give overall accuracy, accuracy per language, micro averaged measures and macro averaged measures. The accuracies are given in percentage and precision, recall and F-measure values are between 0 and 1.

## Program Structure

The structure of the language identification system consists of two parts which are identification models and functions that measures accuracy and metrics of these models. At the beginning of the program, the corpus data is randomly divided into two parts as training and test data such that training data includes 90% of the data while test data includes %10. The training and test sets are kept in list of tuples such that each tuple includes the language and the sentence written in that language.

## Approaches

We used three approaches in our language identification system. The approach to be used can be chosen by giving it as an argument to the program. Note that the space character is omitted in our analysis since we think that it does not give any valuable information about a language.

### Naive Bayes

In the naive bayes model, all sentences of each language are concatenated and acted as one big string. For each character in a language, the probability of being in that language is calculated with laplace smoothing method for the training part. The result is kept in a dictionary whose key is language, and value is a dictionary which includes characters as key and their probabilities as value. In the test part, the probabilities of characters found in the training part in each sentence of test set are added for each language. The result is divided to the number of languages. The language with the highest probability is chosen for that test sentence. The results for each sentences are put in a dictionary whose key is a tuple mentioned above and value is a descending ordered sorted dictionary which contains the probability of being the language of that sentence for each language.

### SVM

We used *SVM<sup>multiclass</sup>* library for SVM learning model. The inputs given that library are prepared in SVM part. Unlike naive bayes, we did not concatenate sentences for each language and determined features and their values for all sentences. Since the library only accepts positive integers as class and feature id, we give unique ids to all languages(classes) and features and put them in dictionaries. After creating input files, for both training and test sets, shell scripts that execute `svm_multiclass_learn` and `svm_multiclass_classify`

are run. The original languages of test sentences and the first line of the output file which includes guessed languages of test sentences are put in a list of tuples.

**Unigram SVM** In the unigram SVM, we selected the existence of characters in each sentence as a feature. We gave a feature id to each character in the corpus, and for each character, if it is in a sentence, that character is given as a feature with value 1 for that sentence.

**Super SVM** In the super SVM, in addition to unigram SVM, we added the number of characters and the number of bigrams in each sentence as features. The bigrams in the corpus are kept in a dictionary with keys as bigrams, and values as ids. The values of character and bigram features are selected as their number of existence for each sentence.

## Performance Measures

We created two functions to measure the performances of the learning models.

### Accuracy

The accuracy function takes tuple list of original and guessed languages as input. It calculates the proportion of all test sentences that are guessed correctly. The accuracy of each language is also calculated and they are returned.

### Metrics

The metrics function takes tuple list of original and guessed languages as input. precision, recall, and F-measure values which are mentioned in evaluation section for both micro and macro averaged cases are calculated. They are returned as values of a dictionary whose keys indicate if they are micro or macro.

## Results

Table 1: Accuracies

	<b>Naive Bayes</b>	<b>unigram SVM</b>	<b>super SVM</b>
<b>Overall Accuracy</b>	<b>32.65384</b>	<b>44.73076</b>	<b>59.30769</b>
Accuracy of bg	0.96618	100.0	67.59776
Accuracy of bs	23.67149	89.78494	26.0
Accuracy of cz	0.0	100.0	98.93617
Accuracy of es-AR	0.0	0.0	1.53846
Accuracy of es-ES	95.83333	0.0	96.78899
Accuracy of hr	27.02702	0.0	55.72139
Accuracy of id	0.94339	95.43378	12.09302
Accuracy of mk	100.0	100.0	100.0
Accuracy of my	100.0	0.0	97.51243
Accuracy of pt-BR	0.0	0.0	91.87817
Accuracy of pt-PT	70.65217	27.52809	16.5
Accuracy of sk	0.0	38.80597	84.81675
Accuracy of sr	1.63043	0.0	23.52941

It is interesting to see that the accuracy for es-ES is around 95%, but es-AR is below 5%. We can say that the system classified the Spanish originated languages to es-ES. We can also clearly say that the languages which have distinct alphabets such as Bulgarian have very high accuracies, which also makes sense since those letters do not appear in any other language. It is also interesting to see that Czech has 0 accuracy in Naive Bayes whereas in SVM, the accuracies are nearly 100%. However, it should be also noted that the sentences are randomly separated into training and test set, so if the test set does not include that language, accuracy would be zero.

It can be also seen that language-specifically, the unigram SVM has accuracies of zero. However, in super SVM, there is never a 0 accuracy, although there are accuracies near 1%.

In addition to these algorithms, we also tried to include number of capital letters as a feature too. The accuracy was better than Naive Bayes and unigram SVM approach; however, it was worse than the super SVM (it was around 50%). Hence, we decided to exclude that feature from our analysis.

Table 2: Metrics

	<b>Naive_Bayes</b>	<b>unigram SVM</b>	<b>super SVM</b>
Micro averaged Precision	0.32653	0.44730	0.59307
Micro averaged Recall	0.32653	0.44730	0.59307
Micro averaged F-measure	0.32653	0.44730	0.59307
Macro averaged Precision	0.29718	0.35139	0.60590
Macro averaged Recall	0.32363	0.42390	0.59454
Macro averaged F-measure	0.20307	0.34042	0.53361

## Improvements and Extensions

In the project, we improved the accuracy of our identification system from approximately 40% to 60% by adding new features to svm input. Although this is a remarkable improvement, there could be more features which will increase the accuracy of the learning model.

We thought about adding word unigrams and average word length as features, but since we eliminated the space characters at the beginning of the program, it would be difficult. Therefore, we decided not to add them.

## Difficulties Encountered

The first difficulty we encountered was because of the project description. In the description, we were told that the corpus is in UTF-8 format. However, we could not process the input properly before encoding it in UTF-16 format. After we discovered that the corpus is in actually UTF-16 format, we could read the file properly.

Implementing our code, we did not face with many issues but we had to choose what we wanted to do in corner cases. For example, when we were implementing averaged measures (precision, recall and F), we got float division by zero error when  $TP_i$  and  $FP_i$  values or  $TP_i$  and  $FN_i$  values are zero at the same time in Naive Bayes. We decided to set the corresponding precision, recall and F values to zero if their denominator is zero.

We also had a "key not found" error when we tried to get the letter probability of a character in a test sentence when that character does not exist in the training set. In this case, instead of adding zero to the probability sum of  $P(c_i|l)$ , we added  $\frac{1}{total\ chars\ in\ l + unique\ chars\ in\ l}$  since we are applying laplace smoothing.

## Conclusion

In our project, we used two models, one discriminative and one generative, to implement a language identification system. We compared these models with some performance measures to see that SVM works better than Naive Bayes. We could improve the accuracies significantly in SVM when we added the count of unigram letters and bigram letters as separate features.



# Bibliography

- [1] [Cornell's SVM-multiclass library]  
[https://www.cs.cornell.edu/people/tj/svm\\_light/svm\\_multiclass.html](https://www.cs.cornell.edu/people/tj/svm_light/svm_multiclass.html)
- [2] [Accuracy, precision, recall, F-measure]  
<http://blog.exsilio.com/all/accuracy-precision-recall-f1-score-interpretation-of-performance-measures/>